

Toward Fast and Accurate Neural Chinese Word Segmentation with Multi-Criteria Learning

Weipeng Huang*, Xingyi Cheng*, Kunlong Chen, Taifeng Wang, Wei Chu
Ant Financial Services Group

{weipeng.hwp|fanyin.cxy|kunlong.ckl|taifeng.wang|weichu.cw}@antfin.com

Abstract

The ambiguous annotation criteria bring into the divergence of Chinese Word Segmentation (CWS) datasets with various granularities. Multi-criteria learning leverage the annotation style of individual datasets and mine their common basic knowledge. In this paper, we proposed a domain adaptive segmenter to capture diverse criteria of datasets. Our model is based on Bidirectional Encoder Representations from Transformers (BERT), which is responsible for introducing external knowledge. We also optimize its computational efficiency via model pruning, quantization, and compiler optimization. Experiments show that our segmenter outperforms the previous results on 10 CWS datasets and is faster than the previous state-of-the-art Bi-LSTM-CRF model.

1 Introduction

Chinese Word Segmentation (CWS) is regarded as a low-level task in NLP. Unlike the language with space between words such as English and French, Chinese is a type of polysynthetic language where compounds are developed from indigenous morphemes (Jernudd and Shapiro, 2011; Gong et al., 2017). The ambiguous distinction between morphemes and compound words leads to the cognitive divergence of words concepts. The labeled datasets seriously diverge due to annotation inconsistency that results in multi-grained compounds. In practice, a segmenter usually provide multiple granularities and configured according to high-level tasks needs. Fine-grained words can help reduce the vocabulary to relieve the sparseness. On the other hand, coarse-grain words make models match exactly and easy to analyze. A multi-criteria model may provide flexibility for this demand.

In recent years, several multi-criteria learning methods of CWS have been proposed to explore the common knowledge of heterogeneous dataset by utilizing the information across the whole corpora, which can boost the out-of-vocabulary (OOV) recalls mutually. (Qiu et al., 2013; Chao et al., 2015; Liu et al., 2016; Chen et al., 2017). First, although the heterogeneous corpora can help each other, the whole datasets are still not big enough to provide adequate linguistic knowledge. Second, the standard recurrent networks including LSTM are limited by decoding speed even using cutting-edge hardware since the computation of states cannot occur in parallel.

In this paper, we propose a multi-criteria method of CWS. Our model uses a domain projection layer to adopt multiple datasets with various granularities. We adopt the bidirectional pre-training encoder from the transformer (BERT) (Vaswani et al., 2017; Devlin et al., 2018) to introduce external knowledge. BERT can be regarded contextual representations and it has achieved great success in some NLU tasks (Reddy et al., 2018; Rajpurkar et al., 2018). But both the fine-tuning and inference procedures of the provided models are computationally inefficient due to a large number of parameters.

The main advantages of our proposed method are scalable and simple, we provide a trade-off between the accuracy and decoding speed. According to the length of the sentence, the number of layers can be adjusted flexibly. We mainly use three techniques including layer-level pruning, quantization and compiler optimization to improve the scalability. Our method not only significantly outperforms the SOTA results on 10 CWS dataset but also faster than the previous SOTA Bi-LSTM-CRF (Ma et al., 2018; Xinchu et al., 2017; Yang et al., 2018) model.

* Equal contribution

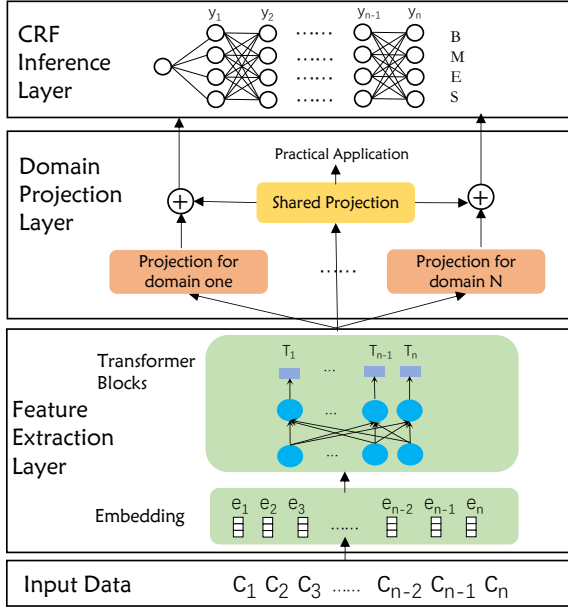


Figure 1: An overview of our model architecture.

2 Model Description

Figure 1 summarizes the proposed model architecture, include a feature extraction layer, a domain projection layer and an inference layer.

2.1 BERT for Feature Extraction

As shown in Figure 1, we employ BERT to extract feature for the input sequence. Characters are first mapped into embedding vectors and then go through several transformer blocks. Compared with Bi-LSTM which process the sequence step by step, the transformer parallelly learns features for all steps so that the decoding speed can be accelerated. However, twelve transformer layers of the original BERT are too heavy for the CWS. To balance computational cost and segmentation accuracy, we prune the layers of BERT and fine tune on our datasets. BERT is pre-trained on a large corpus to capture semantic feature and abundant knowledge, which is of critical importance for the word segmentation task. To speed up both the fine-tuning and inference procedures, we make further optimization as discussed in section 2.4.

2.2 Domain Projection for Multi-Criteria Learning

Inspired by previous works (Chen et al., 2017; Peng and Dredze, 2017), we propose a domain projection layer to enable our model to adapt datasets with diverse criteria. The domain projection layer helps to capture heterogeneous segmen-

tation criteria of each dataset. Section 3.5 shows several examples proving this. There are many variations for the projection layer, while in this paper we use linear transformation which is simple but effective for this task. As shown in Figure 1, an extra shared projection layer is to learn common knowledge from datasets.

2.3 Tag Inference

The output of domain-specific projection and shared projection are concatenated, then feed into the conditional random fields (CRF) layer (Laferty et al., 2001). In CRF layer, the probability of a possible label sequence is formalized as:

$$P(Y|X) = \frac{\prod_{i=2}^n \exp(s(X, i)_{y_i} + b_{y_{i-1}y_i})}{\sum_{y'} \prod_{i=2}^n \exp(s(X, i)_{y'_i} + b_{y'_{i-1}y'_i})} \quad (1)$$

where $y \in \{B, M, E, S\}$ is the label, score function $s(X, i)_{y_i}$ is output of the projection layer at i th character, and $b_{y_{i-1}y_i}$ is trainable parameters. By solving Eq 2 we can obtain the optimal sequence tags:

$$Y^* = \operatorname{argmax} P(Y|X) \quad (2)$$

2.4 Speed Optimization

Neural CWS models improve the performance by increasing the model complexity, which however harms the decoding speed and limits their application in real life. To bridge the gap, we apply model acceleration techniques as follow.

Pruning. Many parameters in deep networks are unimportant or unnecessary, thus pruning methods can be used to remove these parameters and expand the model sparsity (Han et al., 2015). Pruning methods can be categorized into a fine-grained level, kernel level, filter level, and layer level. The acceleration rate increase with the granularity of the pruning strategy. In addition, fine-grained level pruning usually assumes that underlying hardware platforms provide mechanisms for sparse tensor compressing and sparse tensor computation accelerating (Zhu et al., 2017), which is impractical for most current hardware platforms. To maximize the profit of model pruning for real-life application, we perform layer level pruning on the transformer blocks in BERT.

Quantization. Quantization methods also have been investigated for network acceleration. We conduct fixed-point quantization (Gupta et al.,

	PKU	MSR	AS	CITYU	CTB6	SXU	UD	CNC	WTB	ZX
Zhou et al. (2017)	96.0	97.8	-	-	96.2	-	-	-	-	-
Yang et al. (2017)	96.3	97.5	95.7	96.9	96.2	-	-	-	-	-
Chen et al. (2017)	94.3	96.0	94.6	95.6	96.2	96.0	-	-	-	-
Xu and Sun (2017)	96.1	96.3	-	-	95.8	-	-	-	-	-
Yang et al. (2018)	95.9	97.7	-	-	96.3	-	-	-	-	-
Ma et al. (2018)	96.1	97.4	96.2	97.2	96.7	-	96.9	-	-	-
Gong et al. (2018)	96.2	97.8	95.2	96.2	97.3	97.2	-	-	-	-
He (2019)	96.0	97.2	95.4	96.1	96.7	96.4	94.4	97.0	90.4	95.7
Ours (3 layer)	96.6	97.9	96.6	97.6	97.6	97.3	97.3	97.2	93.1	97.0
Ours (3 layer+FP16)	96.5	97.9	96.4	97.5	97.5	97.3	97.3	97.1	92.7	97.0

Table 1: The state of the art performance on different datasets (F-score, %).

	PKU	MSR	AS	CITYU	CTB6	SXU	UD	CNC	WTB	ZX
OOV	3.6	2.7	4.2	7.5	5.6	4.6	12.4	0.7	14.5	5.4
Recall (single-criteria learning)	74.6	78.0	78.1	83.6	61.6	79.4	73.4	64.0	73.3	79.1
Recall (multi-criteria learning)	80.1	84.0	76.9	89.7	88.8	86.0	92.9	62.7	83.7	86.9

Table 2: Test set OOV rate(%), OOV recall(%) achieved with multi-criteria learning, single-criteria learning.

2015) to leverage NVIDIA’s Volta architectural features. Specifically, kernels of multi-head attention layers and feedforward layers use half-precision (FP16), while rest parameters like embedding and normalization parameters use full precision (FP32). The quantization method not only accelerates the computation but also reduce the model size.

Compiler Optimization. XLA is a domain-specific compiler for linear algebra that optimizes TensorFlow (Abadi et al., 2015) computations. By introducing XLA into our model, graphs are compiled into machine instructions, low-level ops are fused to improve the execution speed. For example, **batch matmul** is always followed by a transpose operation in the transformer computation graph. By fusing these two operations, the intermediate product does not need to write back to memory, thus reducing the redundant memory access time and kernel launch overhead.

3 Experiments

3.1 Experimental Settings

Our all experiments are implemented on the hardware with Intel(R) Xeon(R) CPU E5-2682 v4 @ 2.50GHz and NVIDIA Tesla V100.

Datasets. We evaluate our model on ten standard Chinese word segmentation datasets: MSR, PKU, AS, CITYU from SIGHAN 2005 bake-off task (Emerson, 2005). SXU from SIGHAN 2008 bake-off task (MOE, 2008). Chinese Penn Treebank 6.0 (CTB6) from Xue et al. (2005). Chinese Universal Treebank (UD) from the Conll2017 shared task (Zeman and Popel, 2017). WTB (Wang and Yang, 2014), ZX (Zhang and Meishan,

2014) and CNC corpus. For each of the SIGHAN 2005 and 2008 dataset, we randomly select 10% training data as the development set. For other datasets, we use official data split.

Preprocessing. AS and CITYU are mapped from traditional Chinese to simplified Chinese before segmentation. A unique token respectively replaces continuous English characters and digits in the datasets. Full-width tokens are converted to half-width to handle the mismatch between training and test set.

Hyperparameters. The number of domain projection layer is 1, the max sequence length is set to 128. During fine tuning, we use Adam with the learning rate of $2e-5$, L2 weight decay of 0.01, dropout probability of 0.1.

3.2 Main Results

We prune the number of transformer layers from 12 to 1 and find that compared with using 12 layers, the average F-score using 3 layers drop slightly from 97.1% to 96.8% as shown in Table 3. To balance segmentation speed and accuracy, we prune the model to 3 layers. Performance of our model and recent neural CWS models are shown in Table 1. Our model outperform prior work on 10 datasets, with 8.1%, 4.5%, 10.5%, 14.3%, 27.3%, 25.0%, 12.9%, 6.6%, 28.1%, 30.2% error reductions on PKU, MSR, AS, CITYU, CTB6, SXU, UD, CNC, WTB, ZX datasets respectively. Among these datasets, SXU, UD, WTB, ZX are relatively small, but they achieve large error reductions thanks to the shared feature extraction layer. By further applying half-precision (FP16), the accuracy reduction is minor and the model still out-

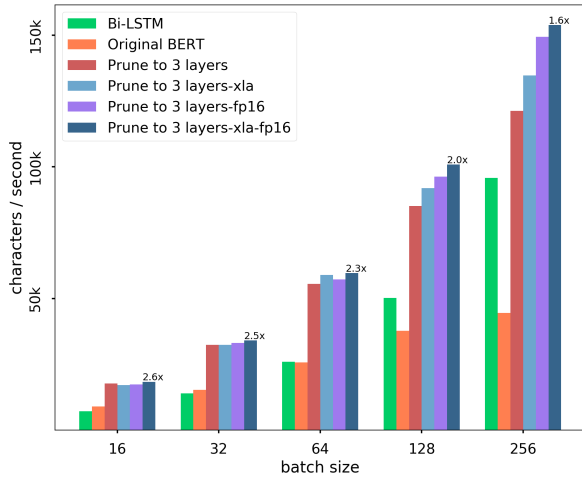


Figure 2: Decoding speed at char level w.r.t batch size, the sequence length is 64.

performs previous SOTA results on 10 datasets.

3.3 Multi-Criteria Learning Improve OOV Recall

Previous work (Huang and Zhao, 2007; Ma et al., 2018) pointed out that OOV is a major error and exploring further sources of knowledge is essential to solving this problem. From a certain point of view, datasets are complementary to each other since OOV in a dataset may appear in other datasets. To utilize knowledge from each other to improve the OOV recall, our model performs multi-criteria learning with the domain projection layer. To evaluate this, we train the proposed model respectively on each dataset, i.e., single-criteria learning. Table 2 shows that comparing with single-criteria learning, multi-criteria learning improve the OOV recall on 10 datasets.

3.4 Scalability

Decoding speed is essential in practice since the word segment is fundamental for many downstream NLP tasks. Previous neural CWS models (Ma et al., 2018; Xinchu et al., 2017; Yang et al., 2018; Gong et al., 2018) use Bi-LSTM with concatenated embedding size 100,100,128,100 respectively. To make a fair comparison, we set the Bi-LSTM embedding size and hidden size to 100, one hidden layer with CRF on the top. Figure 2 shows the decoding speed with regards to batch size. Our model employed original BERT with 12 transformer layers is slower than Bi-LSTM. On the other hand, the speed can be increased by some optimizations including layer-level pruning,

	Precision	Recall	F-Score
Ours (12 layer)	97.2	97.0	97.1
Ours (6 layer)	97.0	97.0	97.0
Ours (3 layer)	96.7	96.9	96.8
Ours (1 layer)	95.6	96.1	95.9

Table 3: Average Precision, Recall, F-score on 10 datasets with different number of layers.

Corpora	PKU,UD,SXU,CITYU,CNC,MSR	AS,CTB6,WTB,ZX	
Segment	副局长	副局长	
Corpora	AS,PKU,CNC,CTB6	CITYU,SXU,UD	MSR
Segment	下午 五时	下午 五 时	下午五时
Corpora	PKU,CTB6,SXU,CITYU	CNC,MSR	AS,UD,ZX,WTB
Segment	令人满意	令人 满意	令人 满意

Figure 3: Our model learns to segment with diverse criteria.

weights quantization, and compiler optimization. Combining all of these three techniques, our models outperform Bi-LSTM with $1.6\times$ - $2.6\times$ acceleration. Our model are more scalable compared with the Bi-LSTM that are limited in their capability to process tasks involving very long sequences. By observing the sequence length distribution, we can search a appropriate layer number to balance F-score and decoding speed.

3.5 Case Study

Figure 3 shows three examples of the segmentation results on all datasets: “下午五时(Five o’clock in the afternoon)”, “副局长(deputy director)”, “令人满意(make sb pleased)”. The segmentation granularity of these words is different according to diverse criteria of the datasets. With the help of the domain projection layer, our model correctly segments these words on each dataset. Without the domain projection layer, the segmentation results are unstable. For instance, “副局长” is segmented as “副局长” or “副/局长” on the same dataset with different context.

4 Conclusion

In this paper, we proposed a simple but effective Chinese Word Segmentation (CWS) method that employ BERT and add a domain projection layer on the top with multi-criteria learning. To be practicability, acceleration techniques including pruning, quantization, and compiler optimization are applied to improve the word segmentation speed. Experiments show that our proposed model achieve higher performance on the CWS accuracy and prediction speed than the SOTA methods.

References

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. [TensorFlow: Large-scale machine learning on heterogeneous systems](#). Software available from tensorflow.org.
- Jiayuan Chao, Zhenghua Li, Wenliang Chen, and Min Zhang. 2015. Exploiting heterogeneous annotations for weibo word segmentation and pos tagging. In *Natural Language Processing and Chinese Computing*, pages 495–506. Springer.
- X. Chen, Z. Shi, X. Qiu, and X. Huang. 2017. Adversarial multi-criteria learning for chinese word segmentation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, volume 1, pages 1193–1203.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Thomas Emerson. 2005. The second international chinese word segmentation bakeoff. In *Proceedings of the fourth SIGHAN workshop on Chinese language Processing*.
- Chen Gong, Zhenghua Li, Min Zhang, and Xinzhou Jiang. 2017. Multi-grained chinese word segmentation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 692–703.
- Jingjing Gong, Xinchu Chen, Tao Gui, and Xipeng Qiu. 2018. Switch-lstms for multi-criteria chinese word segmentation. *arXiv preprint arXiv:1812.08033*.
- Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. 2015. Deep learning with limited numerical precision. In *International Conference on Machine Learning*, pages 1737–1746.
- Song Han, Huizi Mao, and William J Dally. 2015. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*.
- Han He. 2019. Effective neural solution for multi-criteria word segmentation. pages 133–142. *Smart Intelligent Computing and Applications*. Springer.
- Changning Huang and Hai Zhao. 2007. Chinese word segmentation: A decade review. *Journal of Chinese Information Processing*, 21(3):8–20.
- Björn H Jernudd and Michael J Shapiro. 2011. *The politics of language purism*, volume 54. Walter de Gruyter.
- J Lafferty, A McCallum, and F C N Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data.
- Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. 2016. Deep multi-task learning with shared memory. *arXiv preprint arXiv:1609.07222*.
- Ji Ma, Kuzman Ganchev, and David Weiss. 2018. State-of-the-art chinese word segmentation with bi-lstms. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- PRC MOE. 2008. The fourth international chinese language processing bakeoff: Chinese word segmentation, named entity recognition and chinese pos tagging. In *Proceedings of the sixth SIGHAN workshop on Chinese language processing*.
- N Peng and M Dredze. 2017. Multi-task domain adaptation for sequence tagging. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pages 91–100.
- Xipeng Qiu, Jiayi Zhao, and Xuanjing Huang. 2013. Joint chinese word segmentation and pos tagging on heterogeneous annotated corpora with multiple task learning. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 658–668.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for squad. *arXiv preprint arXiv:1806.03822*.
- Siva Reddy, Danqi Chen, and Christopher D Manning. 2018. Coqa: A conversational question answering challenge. *arXiv preprint arXiv:1808.07042*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.
- Wang and William Yang. 2014. Dependency parsing for weibo: An efficient probabilistic logic programming approach. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1152–1158.
- Chen Xinchu, Zhan Shi, Xipeng Qiu, and Xuanjing Huang. 2017. Adversarial multi-criteria learning for chinese word segmentation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, volume 1, pages 1193–1203. Association for Computational Linguistics.

J Xu and X Sun. 2017. Dependency-based gated recursive neural network for chinese word segmentation. In *The 54th Annual Meeting of the Association for Computational Linguistics*, pages 1193–1203.

Naiwen Xue, Fei Xia, Fu-Dong Chiou, and Marta Palmer. 2005. The penn chinese treebank: Phrase structure annotation of a large corpus. *Natural language engineering*, 11(02):207–238.

J Yang, Y Zhang, and F Dong. 2017. Neural word segmentation with rich pretraining. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 839–849.

J Yang, Y Zhang, and S Liang. 2018. Subword encoding in lattice lstm for chinese word segmentation. In *arXiv preprint arXiv:1810.12594*.

Zeman and Martin Popel. 2017. Conll 2017 shared task: Multilingual parsing from raw text to universal dependencies. In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 1–19. Association for Computational Linguistics.

Zhang and Meishan. 2014. Type-supervised domain adaptation for joint segmentation and pos-tagging. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 588–597.

H. Zhou, Z. Yu, Y. Zhang, S. Huang, X.-Y. DAI, and J. Chen. 2017. Word-context character embeddings for chinese word segmentation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 133–142.

Michael Zhu, Suyog Gupta, Michael Zhu, Suyog Gupta, Michael Zhu, and Suyog Gupta. 2017. To prune, or not to prune: exploring the efficacy of pruning for model compression.

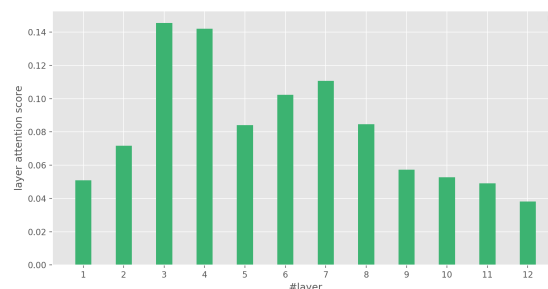


Figure 4: Distribution of layer attention score.

A Do we really need 12 transformer layers for word segmentation?

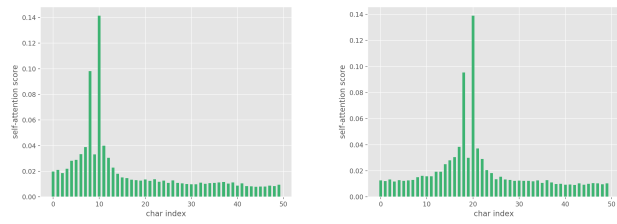
In our paper, We have demonstrated that BERT with multi-criteria learning has superior performance on accuracy and prediction speed in word segmentation. We use 3 transformer layers instead of the original 12 layers to balance the F-score and decoding speed. Why do three layers seem to the CWS on ten datasets to be the most cost-effective? Here we do some analysis as complementary.

A.1 Layer Attention

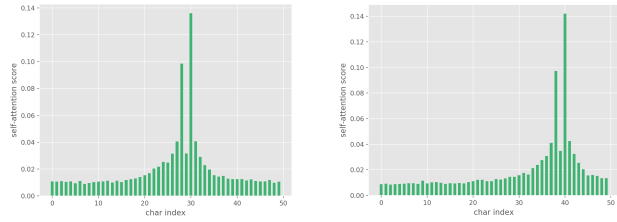
BERT has achieved great success in many NLU tasks by pre-training a stack of 12 transformer layers to learn abundant knowledge. Intuitively, top layers capture high-level semantic features while bottom layers learn low-level features like grammar. As for word segmentation task, high-level semantic features may have a small impact (see the section A.2) so that we make further investigation to find the minimal number of transformer layers. We freeze weights of each layer in the pre-trained BERT and conduct layer attention fine-tuning on word segmentation datasets. As shown in Figure 4, the attention score gradually decrease in top layers from 7 to 12, and **the third** layer gains the highest attention score. The results prove that the model with three layers contains most information for word segmentation.

A.2 Self Attention

In the self-attention layers of the transformer, the attention score of each character is calculated with the rest characters. We average the attention scores at each index of the sentences with a length larger than 50. As shown in Figure 5, characters around the current character gain larger weights than those far away. The result indicates that word segmentation depends more on grammar and long term dependencies are relatively unimportant. It



(a) Current char index 10 (b) Current char index 20



(c) Current char index 30 (d) Current char index 40

Figure 5: Distribution of self-attention score at each char index, char index 10, 20, 30, 40 are shown.

intuitively proves that it is not necessary to keep long term memory of the sequence for CWS. As for the transformer, self-attention may be limited to a fixed window size to reduce computation and make model acceleration. We leave this for future work.